

Code Generation: Failure of Judgment Analysis

Table 1: **(Java)** False Negatives: analysis of the causes of the misjudgment of the LLMs. Every cell repots the number of occurrences of each category on a sample of 15 cases of misjudgment. Note that more than one category can be assigned to each assignment.

	DSC 1.3B	DSC 6.7B	DSC 33B	CL 7B	CL 13B	CL 34B	GPT 3.5	GPT 4
<i>Test result unreliable</i>	1	0	3	0	0	1	3	6
<i>Ambiguous docstring</i>	0	0	1	1	0	0	2	0
<i>Artificial hallucination</i>	7	0	6	11	9	7	5	0
<i>Uncought wrong behavioral</i>	0	0	0	0	0	0	0	0
<i>Misunderstanding of code statements</i>	7	0	3	4	6	7	3	1
<i>Focus on non-functional requirements</i>	0	0	3	0	0	1	0	0
<i>Limited coding context</i>	1	0	0	0	0	2	7	10
<i>Misintepreted implementation requirements</i>	0	0	2	1	0	1	0	1
<i>Shallow description</i>	0	0	0	0	0	0	0	0

Table 2: **(Java)** False Positives: analysis of the causes of the misjudgment of the LLMs. Every cell repots the number of occurrences of each category on a sample of 15 cases of misjudgment. Note that more than one category can be assigned to each assignment.

	DSC 1.3B	DSC 6.7B	DSC 33B	CL 7B	CL 13B	CL 34B	GPT 3.5	GPT 4
<i>Test result unreliable</i>	1	0	0	0	0	0	0	1
<i>Ambiguous docstring</i>	4	4	3	1	4	3	4	2
<i>Artificial hallucination</i>	0	0	0	0	0	0	0	0
<i>Uncought wrong behavioral</i>	6	3	10	12	10	7	7	9
<i>Misunderstanding of code statements</i>	0	0	0	0	0	0	0	0
<i>Focus on non-functional requirements</i>	0	0	0	0	0	0	0	0
<i>Limited coding context</i>	9	11	6	1	4	8	7	8
<i>Misintepreted implementation requirements</i>	0	1	0	0	0	0	0	0
<i>Shallow description</i>	0	0	0	0	0	0	0	0

Table 3: **(Python)** False Negatives: analysis of the causes of the misjudgment of the LLMs. Every cell repots the number of occurrences of each category on a sample of 15 cases of misjudgment. Note that more than one category can be assigned to each assignment.

	DSC 1.3B	DSC 6.7B	DSC 33B	CL 7B	CL 13B	CL 34B	GPT 3.5	GPT 4
<i>Test result unreliable</i>	3	2	1	1	4	3	3	2
<i>Ambiguous docstring</i>	0	1	2	0	1	2	1	5
<i>Artificial hallucination</i>	1	8	6	4	6	3	4	1
<i>Uncought wrong behavioral</i>	0	0	0	0	0	0	0	0
<i>Misunderstanding of code statements</i>	3	2	0	1	3	2	1	3
<i>Focus on non-functional requirements</i>	0	1	1	0	0	0	0	1
<i>Limited coding context</i>	1	1	2	0	0	1	5	6
<i>Misintepreted implementation requirements</i>	4	0	1	1	0	1	2	1
<i>Shallow description</i>	1	0	1	8	3	3	0	0

Table 4: **(Python)** False Positives: analysis of the causes of the misjudgment of the LLMs. Every cell repots the number of occurrences of each category on a sample of 15 cases of misjudgment. Note that more than one category can be assigned to each assignment.

	DSC 1.3B	DSC 6.7B	DSC 33B	CL 7B	CL 13B	CL 34B	GPT 3.5	GPT 4
<i>Test result unreliable</i>	0	0	0	0	0	0	0	0
<i>Ambiguous docstring</i>	7	7	5	7	4	2	6	6
<i>Artificial hallucination</i>	0	0	0	0	0	1	0	0
<i>Uncought wrong behavioral</i>	2	6	4	4	4	8	5	2
<i>Misunderstanding of code statements</i>	0	0	0	0	0	0	0	0
<i>Focus on non-functional requirements</i>	0	0	0	0	0	0	0	0
<i>Limited coding context</i>	5	2	9	3	1	1	4	7
<i>Misintepreted implementation requirements</i>	0	1	0	0	0	0	0	0
<i>Shallow description</i>	1	0	0	2	6	4	0	0